

1 **Iroki: automatic customization for phylogenetic trees**

2

3 Ryan M. Moore^{1*}, Amelia O. Harrison², Sean M. McAllister⁴, Rachel Marine³, Clara Chan⁴,
4 and K. Eric Wommack¹

5

6 ¹Center for Bioinformatics & Computational Biology, University of Delaware, Newark, DE,
7 USA

8 ²Department of Entomology and Wildlife Ecology, University of Delaware, Newark, DE,
9 USA

10 ³Centers for Disease Control, Atlanta, Georgia, USA

11 ⁴School of Marine Science and Policy, University of Delaware, Newark, DE, USA

12

13 Corresponding author's information

14 *To whom correspondence should be addressed

15 **Address: Delaware Biotechnology Institute, 15 Innovation Way, Newark, Delaware**

16 **19711**

17 **(Tel): (302) 831-4362**

18 **(Fax): (302) 831-3447**

19 **(E-mail): moorer@udel.edu**

20

21 **Abstract**

22 *Background*

23 Phylogenetic trees are an important analytical tool for examining species and community
24 diversity and the evolutionary history of species. In the case of microorganisms,
25 decreasing sequencing costs have enabled researchers to generate ever-larger sequence
26 datasets, which in turn have begun to fill gaps in the evolutionary history of microbial
27 groups. However, phylogenetic analyses of large sequence datasets present challenges to
28 extracting meaningful trends from complex trees. Scientific inferences made by visual
29 inspection of phylogenetic trees can be simplified and enhanced by customizing various
30 parts of the tree, including label color, branch color, and other features. Yet, manual
31 customization is time-consuming and error prone, and programs designed to assist in
32 batch tree customization often require programming experience. To address these
33 limitations, we developed Iroki, a program for fast, automatic customization of
34 phylogenetic trees. Iroki allows the user to incorporate information on a broad range of
35 metadata for each experimental unit represented in the tree.

36

37 *Results*

38 Iroki was applied to four existing microbial sequence datasets to demonstrate its utility in
39 data exploration and presentation. Iroki was used to highlight connections between viral
40 phylogeny and host taxonomy, to explore the abundance of microbial groups across
41 samples of cattle hide, to examine short-term temporal dynamics of viroplankton
42 communities, and to search for trends in the biogeography of Zetaproteobacteria.

43

44 *Conclusions*

45 Iroki is an easy-to-use web app and command line application for fast, automatic
46 customization of phylogenetic trees based on user-provided categorical or continuous
47 metadata. Iroki allows for rapid hypothesis testing through visualizing custom
48 phylogenetic trees, streamlining the process of phylogenetic data exploration and
49 presentation.

50

51 **Availability**

52 Iroki can be accessed through a web app or via installation through RubyGems, from
53 source, or through the Iroki Docker image. All source code and documentation is
54 available under the GPLv3 license at <https://github.com/mooreryan/iroki>. The Iroki web-
55 app is accessible at www.iroki.net or through the Virome portal
56 (<http://virome.dbi.udel.edu>), and its source code is released under GPLv3 license at
57 https://github.com/mooreryan/iroki_web. The Docker image can be found here:
58 <https://hub.docker.com/r/mooreryan/iroki>.

59

60 **Keywords**

61 Phylogeny, visualization, sequence analysis, bioinformatics, metagenomics

62

63 **Iroki: automatic customization for phylogenetic trees**

64

65 **Background**

66 Studies in microbial ecology often use phylogenetic trees as a means for assessing the
67 diversity and evolutionary history of microorganisms. As the cost of sequencing has
68 declined, researchers have been able to gather ever-larger sequence datasets. While large
69 sequence datasets have begun to fill in the gaps in the evolutionary history of microbial
70 groups [1–5]; they have also posed new analytical challenges as extracting meaningful
71 trends within such highly dimensional datasets can be cumbersome. In particular,
72 scientific inferences made by visual inspection of phylogenetic trees can be simplified and
73 enhanced by customizing various parts of the tree including label and branch color,
74 branch width, and other features. Though many tree visualization packages allow for
75 manual modifications [6–9], the process can be time consuming and error prone
76 especially when the tree contains many nodes. While a handful of existing programs
77 address the issue of tree visualization, most are not capable of batch customization and
78 those that do often require programming experience [10–13].

79

80 Iroki, a program for fast, automatic customization of phylogenetic trees, was developed to
81 address these limitations and enable users to incorporate information on a broad range of
82 metadata for each experimental unit represented in the tree. Iroki is available for use
83 through a web interface at www.iroki.net, through the Virome portal
84 (<http://virome.dbi.udel.edu>), and through a UNIX command line tool. Results are saved in

85 the widely used Nexus format with color metadata tailored for use with FigTree [8] (a
86 freely available and efficient tree viewer).

87

88 **Implementation**

89 Iroki enhances visualization of phylogenetic trees by coloring node labels and branches
90 according to categorical metadata criteria or numerical data such as abundance
91 information. Iroki can also rename nodes in a batch process according to user
92 specifications so that node names are more descriptive. A tree file in Newick format
93 containing a phylogenetic tree is always required. Additional required input files depend
94 on the operation(s) desired. Coloring functions require a color map or a biom [14] file.
95 Node renaming functions require a name map. The color map, name map, and biom files
96 are created by the user and, along with the Newick file, form the inputs for Iroki.

97

98 *Explicit tree coloring*

99 Iroki's principle functionality involves coloring node labels and/or branches based on
100 information provided by the user in the color map. The color map text file contains either
101 two or three tab-delimited columns depending on how branches and labels are to be
102 colored. Two columns, pattern and color, are used when labels and branches are to have
103 the same color. Three columns, pattern, label color, and branch color, are used when
104 branches and labels are to have different colors. Patterns are searched against node labels
105 either as regular expressions or exact string matches.

106

107 Entries in the color column can be any of the 657 named colors in the R programming
108 language [15] (e.g., skyblue, tomato, goldenrod2, lightgray, black) or any valid
109 hexadecimal color code (e.g., #FF78F6). In addition, Iroki provides a 19 color palette
110 with complementary colors based on Kelly's color scheme for maximum contrast [16].
111 Nodes in the tree that are not in the color map will remain black.

112

113 Depending on user-specified options, a pattern match to node label(s) will trigger coloring
114 of the label and/or the branch directly connected to that label. Inner branches will be
115 colored to match their descendent branches if all descendants are the same color,
116 allowing quick identification of common ancestors and clades that share common
117 metadata.

118

119 *Tree coloring based on numerical data*

120 Iroki provides the ability to generate color gradients based on numerical data, such as
121 absolute or relative abundance, from a tab-delimited biom format file. Single-color
122 gradients use color saturation to illustrate numerical differences, with nodes at a higher
123 level being more saturated than those at a lower level. For example, highly abundant
124 nodes will be represented by more highly saturated colors. Two-color gradients show
125 numerical differences through both color mixing and luminosity. Additionally, the biom
126 file may specify numerical information for one group (e.g., abundance in a particular
127 sample) or for two groups (e.g., abundance in the treatment group vs. abundance in the

128 control group). For biom files with one group, single- or two-color gradients may be used.
129 However, biom files specifying two-group metadata may only use the two-color gradient.

130

131 *Renaming nodes*

132 Some packages for generating phylogenetic trees (RAxML [17], PHYLIP [18], etc.) require
133 node names to be ten characters or less. Name restrictions present challenges to scientific
134 interpretation of phylogenetic trees. Iroki's renaming function uses a two-column, tab-
135 delimited name map to associate current node names, exactly matching those in the tree
136 file, with new names. The new name column has no restrictions on name length or
137 character type. Iroki ensures name uniqueness by appending integers to the ends of
138 names, if necessary.

139

140 *Combining the color map, name map, and biom files*

141 Iroki can be used to make complex combinations of customizations by combining the
142 color map, name map, and biom files. For example, a biom file can be used to apply a
143 color gradient based on numerical data to the labels of a tree, a color map can be used to
144 separately color the branches based on user-specified conditions, and a name map can be
145 used to rename nodes in a single command or web request. Iroki follows a specific order
146 of precedence when applying multiple customizations. First, the color gradient inferred
147 from the biom file is applied. Next, the color map is applied to specified labels or
148 branches, overriding the gradient applied in the previous step if necessary. Finally, the
149 name map is used to map current names to the new names (Fig. 1).

150

151 *Output*

152 Iroki outputs the modified tree in the Nexus format. When building the phylogenetic tree,

153 FigTree uses the Nexus format file and interprets the color metadata output of Iroki.

154

155 **Results & Discussion**

156

157 *Global diversity of bacteriophage*

158 Viruses are the most abundant biological entity on Earth, providing an enormous reservoir

159 of genetic diversity, driving evolution of their hosts, influencing composition of microbial

160 communities, and affecting global biogeochemical cycles [19,20]. The current viral

161 taxonomic system is based on a suite of physical characteristics of the virion rather than

162 on genome sequences. The phage proteomic tree, created to provide a genome-based

163 taxonomic system for bacteriophage classification [21], was recently updated to include

164 hundreds of new phage genomes from the Phage SEED reference database [22], as well as

165 long assembled contigs from viral shotgun metagenomes (viromes) collected from the

166 Chesapeake Bay (SERC) [23] and the Mediterranean Sea [24].

167

168 Taxonomy and host information metadata was collected for the viral genome sequences, a

169 color map was created to assign colors based on viral family and host phyla, and Iroki was

170 used to add color metadata to branches and labels of the phage proteomic tree. Since a

171 large number of colors were required on the tree, Iroki's Kelly color palette was used to
172 provide clear color contrasts. The tree was rendered with FigTree (Fig. 2).

173

174 Adding color to the phage proteomic tree with Iroki shows trends in the data that would
175 be difficult to discern without color. Uncultured phage contigs from the SERC and
176 Mediterranean viromes make up a large portion of all phage sequences shown on the tree,
177 and are widely distributed among known phage. In general, viruses in the same family
178 claded together, e.g., branch coloring highlights large groups of closely related
179 Siphoviridae and Myoviridae. This label-coloring scheme also shows that viruses infecting
180 hosts within same phylum are, in general, phylogenetically similar. For example, viruses
181 within one of the multiple large groups of Siphoviridae across the tree infect almost
182 exclusively host species within the same phylum, e.g., Siphoviridae infecting
183 Actinobacteria clade away from Siphoviridae infecting Firmicutes or Proteobacteria.

184

185 *Bacterial community diversity and prevalence of E. coli in beef cattle*

186 Shiga toxin-producing *Escherichia coli* (STEC) are dangerous human pathogens that
187 colonize the lower gastrointestinal tracts of cattle and other ruminants. STEC-
188 contaminated beef and STEC shed in the feces of these animals are major sources of
189 foodborne illness. To identify possible interactions between STEC populations and the
190 commensal cattle microbiome, a recent study examined diversity of the bacterial
191 community associated with beef cattle hide [25]. Fecal and hide samples were collected
192 over twelve weeks and SSU rRNA amplicon libraries were constructed and analyzed by

193 Illumina sequencing [26]. The study indicated that the community structure of hide
194 bacterial communities was altered when the hides were positive for STEC contamination.
195
196 Iroki was used to visualize changes in the relative abundance of each cattle hide bacterial
197 OTU according to the presence or absence of STEC. A Mann-Whitney U test comparing
198 OTU abundance between STEC positive and STEC negative samples was performed, and
199 those bacterial OTUs showing a significant change in relative abundance ($p < 0.5$) were
200 placed on a phylogenetic tree according to the 16S rRNA sequence. Branches of the tree
201 were colored based on whether there was a significant change in relative abundance with
202 STEC contamination (red: $p \leq 0.05$, blue: $p > 0.05$). Node labels were colored along a
203 blue-green color gradient representing the abundance ratio of OTUs between samples
204 with STEC (blue) and without (green). Additionally, label luminosity was determined
205 based on overall abundance (lighter: less abundant, darker: more abundant) (Fig. 3). Iroki
206 makes it clear that most OTUs on the tree showed a significant difference in abundance
207 (branch coloring) between STEC positive and STEC negative samples (node coloring).
208 Furthermore, we can see that most OTUs are at low abundance with only a few highly
209 abundant OTUs (label luminosity). The color gradient added by Iroki allows us to see that
210 the abundant OTUs were evolutionarily distant from one another and thus spread out
211 across many phylogenetic groups.
212
213 Iroki can be used to quickly test hypotheses without investing a large amount of time
214 annotating trees manually. A UPGMA tree was created based on unweighted UniFrac

215 distance [27] between 356 bacterial community profiles based on SSU rRNA amplicon
216 sequences from cattle hide and feces samples (Fig. 4). Iroki was used to evaluate
217 similarities in sample bacterial communities according to the sampling location. Iroki
218 colored branches based on whether the sample originated from feces (blue) or from hide
219 (red). The coloring added by Iroki shows a clear partitioning of bacterial communities on
220 the tree based on their sampling location (hide or feces). However, four fecal samples
221 claded with hide samples, and two hide samples claded with fecal samples, making them
222 be good candidates for more in-depth examination. Additionally, Iroki was used to
223 illustrate a correlation between one of the most abundant bacterial families,
224 Ruminococcaceae, and sampling location. Iroki colored node labels with a color gradient
225 based on Ruminococcaceae family abundance, utilizing both a single color gradient (Fig.
226 4A) and a two color gradient (Fig. 4B). Custom trees were visualized using FigTree.
227 Iroki's automatic color gradient and ability to label branches and nodes based on different
228 criteria clearly show that Ruminococcaceae is more abundant in fecal samples than in
229 hide samples.

230

231 *Short-term dynamics of virioplankton*

232 The gene encoding Ribonucleotide reductase (RNR) is common within viral genomes and
233 thus can be used as a marker gene for studying viral diversity. Moreover, RNR
234 polymorphism is predictive of some of the biological and ecological features of viral
235 populations [28]. A mesocosm experiment examined the short-term dynamics of phage
236 populations using RNR amplicon sequences, specifically, sequences of class II RNRs of

237 bacteriophages infecting cyanobacterial hosts. A phylogenetic tree was created from the
238 Cyano II RNR amplicon sequences and Iroki was used to color nodes and branches based
239 on the time point (0 h, 6 h, 12 h) at which each amplicon sequence was observed. The
240 customized tree was then visualized using FigTree (Fig. 5). Iroki's coloring shows that no
241 phylogenetic clade was dominated by OTUs observed in any particular time point; rather,
242 time points were spread relatively evenly across clades. This analysis demonstrates Iroki's
243 utility for exploring sequence datasets, allowing the researcher to quickly and easily test
244 hypotheses.

245

246 *Phylogeny of Zetaproteobacteria within a biogeographic context*

247 Biogeographical studies assess the distribution of an organism's biodiversity across space
248 and time. The extent to which microorganisms exhibit biogeography is an open question
249 in microbial ecology. The isolated nature of the microbial communities associated with
250 deep-ocean hydrothermal vents, provides an ideal system for studying the biogeography of
251 microbes. In particular, iron-oxidizing bacteria have been shown to thrive in vent fluids,
252 sediments, and iron-rich microbial mats associated with the vents. Globally, iron-
253 oxidizing bacteria make significant contributions to the iron and carbon cycles. A recent
254 study analyzed multiple SSU rRNA clone libraries to investigate the biogeography of
255 Zetaproteobacteria, a phyla containing many iron-oxidizing bacterial species, between
256 three sampling regions of the Pacific Ocean (central Pacific—Loihi seamount, western
257 Pacific—Southern Mariana Trough, and southern Pacific (Vailulu'u Seamount/Tonga
258 Arc/East Lau Spreading Center/Kermadec Arc) [29]. Sequences were aligned and a

259 phylogeny was inferred as described in [29]. Iroki was used to examine the relationship
260 between sampling location and phylotype by adding branch and label color based on
261 geographic location and renaming original node labels with OTU and location metadata.
262 The custom tree was visualized using FigTree (Fig. 6). In some cases, OTUs contained
263 sequences from only one sampling location (e.g., OTUs 12, 15, and 16), whereas other
264 OTUs are distributed among more than one sampling location (e.g., OTUs 1, 2, and 4).
265 Often, sequences samples from the same geographic location are in the same phylotype
266 despite being members of different OTUs (e.g., OTUs 10 and 19).

267

268 *Availability and requirements*

269 A web-based version of Iroki can be accessed online at www.iroki.net or through the
270 Virome portal (<http://virome.dbi.udel.edu/>). For users who wish to run Iroki locally, a
271 command line version of the program is installable via RubyGems, from GitHub
272 (<https://github.com/mooreryan/iroki>). A Docker image is available for users who desire the
273 flexibility of the command line tool, but do not want to install Iroki or manage its
274 dependencies (<https://hub.docker.com/r/mooreryan/iroki>). Docker is a convenient method
275 for packaging an application with all of its dependencies that is guaranteed to run the
276 same regardless of the user's environment [30,31]. The README provided with the source
277 code provides detailed instructions for setting up and running Iroki. Further
278 documentation and tutorials can be found at the Iroki Wiki
279 (<https://github.com/mooreryan/iroki/wiki>).

280

281 *License*

282 Iroki and its associated programs are released under the GNU General Public License
283 version 3 [32].

284

285 **Conclusions**

286 Iroki is a command line program and web app for fast, automatic customization of large
287 phylogenetic trees based on user specified configuration files describing categorical or
288 continuous metadata information. The output files include Nexus tree files with color
289 metadata tailored specifically for use with FigTree. Various example datasets from
290 microbial ecology studies were analyzed to demonstrate Iroki's utility. In each case, Iroki
291 simplified the processes of data exploration, data presentation, and hypothesis testing.
292 Iroki provides a simple and convenient way to rapidly customize phylogenetic trees,
293 especially in cases where the tree in question is too large to annotate manually or in
294 studies with many trees to annotate.

295

296 **List of Abbreviations**

297 OTU: operational taxonomic

298 RNR: Ribonucleotide reductase

299 STEC: Shiga-toxigenic *Escherichia coli*

300

301 **Ethics approval and consent to participate**

302 Not applicable

303

304 **Consent for publication**

305 Not applicable

306

307 **Availability of data and materials**

308 Data and code used to generate figures are available on GitHub at

309 https://github.com/mooreryan/iroki_manuscript_data

310

311 **Funding**

312 This project was supported by the USDA National Institute of Food and Agriculture award

313 number 2012-68003-30155 and the National Science Foundation Advances in

314 Bioinformatics program (award number DBI_1356374).

315

316 **Competing Interests**

317 The authors declare that they have no competing interests.

318

319 **Authors' contributions**

320 RMM and SMM conceived the project. RMM wrote the manuscript and implemented

321 Iroki. AOH and RM processed and analyzed Cyano II amplicons. All authors read,

322 edited, and approved the final manuscript.

323

324 **Acknowledgements**

325 We would like to acknowledge Daniel J. Nasko and Jessica M. Chopyk for their work on
326 the phage proteomic tree, and Barbra D. Ferrell for editing the manuscript.
327

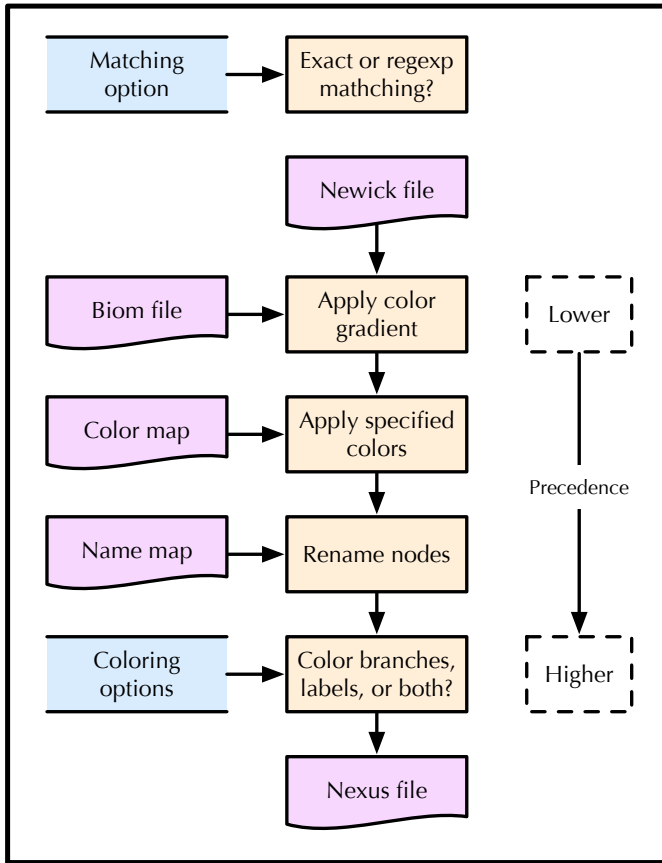
328 **References**

- 329 1. Lan Y, Rosen G, Hershberg R. Marker genes that are less conserved in their sequences
330 are useful for predicting genome-wide similarity levels between closely related prokaryotic
331 strains. *Microbiome*. 2016;4:18.
- 332 2. Larkin A a, Blinebry SK, Howes C, Lin Y, Loftus SE, Schmaus CA, et al. Niche
333 partitioning and biogeography of high light adapted *Prochlorococcus* across taxonomic
334 ranks in the North Pacific. *ISME J*. 2016;1–13.
- 335 3. Simister RL, Deines P, Botté ES, Webster NS, Taylor MW. Sponge-specific clusters
336 revisited: A comprehensive phylogeny of sponge-associated microorganisms. *Environ.*
337 *Microbiol*. 2012;14:517–24.
- 338 4. Wu Z, Yang L, Ren X, He G, Zhang J, Yang J, et al. Deciphering the bat virome catalog
339 to better understand the ecological diversity of bat viruses and the bat origin of emerging
340 infectious diseases. *ISME J*. 2016;10:609–20.
- 341 5. Müller AL, Kjeldsen KU, Rattei T, Pester M, Loy A. Phylogenetic and environmental
342 diversity of DsrAB-type dissimilatory (bi)sulfite reductases. *ISME J*. 2015;9:1152–65.
- 343 6. University W. Phylogeny Programs.
344 <http://evolution.genetics.washington.edu/phylip/software.html#Plotting>. Accessed 2016 Jul
345 21.
- 346 7. Zhang H, Gao S, Lercher MJ, Hu S, Chen WH. EvolView, an online tool for visualizing,
347 annotating and managing phylogenetic trees. *Nucleic Acids Res*. 2012;40.
- 348 8. Rambaut A. FigTree. <http://tree.bio.ed.ac.uk/software/figtree/>. Accessed 2016 Jul 21.
- 349 9. Zmasek CM. Archaeopteryx.

- 350 <https://sites.google.com/site/cmzmasek/home/software/archaeopteryx>. Accessed 2016 Jul
351 21.
- 352 10. Paradis E, Claude J, Strimmer K. APE: Analyses of phylogenetics and evolution in R
353 language. *Bioinformatics*. 2004;20:289–90.
- 354 11. Revell LJ. phytools: An R package for phylogenetic comparative biology (and other
355 things). *Methods Ecol. Evol.* 2012;3:217–23.
- 356 12. Huerta-Cepas J, Dopazo J, Gabaldón T. ETE: a python Environment for Tree
357 Exploration. *BMC Bioinformatics*. 2010;11:24.
- 358 13. Chen W-H, Lercher MJ, Ganfornina M, Gutierrez G, Bastiani M, Sanchez D, et al.
359 ColorTree: a batch customization tool for phylogenetic trees. *BMC Res. Notes. BioMed*
360 *Central*; 2009;2:155.
- 361 14. McDonald D, Clemente JC, Kuczynski J, Rideout J, Stombaugh J, Wendel D, et al. The
362 Biological Observation Matrix (BIOM) format or: how I learned to stop worrying and love
363 the ome-ome. *Gigascience*. 2012;1:7.
- 364 15. Ripley BD. The R project for statistical computing. 2001. p. 1–3.
- 365 16. Kelly KL. Twenty-two colors of maximum contrast. *Color Eng.* 1965. p. 26–7.
- 366 17. Stamatakis A. RAxML version 8: A tool for phylogenetic analysis and post-analysis of
367 large phylogenies. *Bioinformatics*. 2014;30:1312–3.
- 368 18. Felsenstein J. PHYLIP. <http://evolution.gs.washington.edu/phytip.html>. Accessed 2016
369 Jul 21.
- 370 19. Suttle CA. Marine viruses — major players in the global ecosystem. *Nat. Rev.*
371 *Microbiol.* Nature Publishing Group; 2007;5:801–12.

- 372 20. Rohwer F, Thurber RV. Viruses manipulate the marine environment. *Nature*. Nature
373 Publishing Group; 2009;459:207–12.
- 374 21. Rohwer F, Edwards R. The Phage Proteomic Tree: a genome-based taxonomy for
375 phage. *J. Bacteriol.* 2002;184:4529–35.
- 376 22. Phage SEED. <http://www.phantome.org/PhageSeed/Phage.cgi>. Accessed 2016 Jul 21.
- 377 23. Wommack KE, Nasko DJ, Chopyk J, Sakowski EG. Counts and sequences,
378 observations that continue to change our understanding of viruses in nature. *J. Microbiol.*
379 2015;53:181–92.
- 380 24. Mizuno CM, Rodriguez-Valera F, Kimes NE, Ghai R. Expanding the marine virosphere
381 using metagenomics. *PLoS Genet. Public Library of Science*; 2013;9:e1003987.
- 382 25. Chopyk J, Moore RM, DiSpirito Z, Stromberg ZR, Lewis GL, Renter DG, et al. Presence
383 of pathogenic *Escherichia coli* is correlated with bacterial community diversity and
384 composition on pre-harvest cattle hides. *Microbiome. BioMed Central*; 2016;4:9.
- 385 26. Fadrosch DW, Ma B, Gajer P, Sengamalay N, Ott S, Brotman RM, et al. An improved
386 dual-indexing approach for multiplexed 16S rRNA gene sequencing on the Illumina
387 MiSeq platform. *Microbiome*. 2014;2:6.
- 388 27. Lozupone C, Knight R. UniFrac: a New Phylogenetic Method for Comparing Microbial
389 Communities. *Appl. Environ. Microbiol. American Society for Microbiology*;
390 2005;71:8228–35.
- 391 28. Sakowski EG, Munsell E V., Hyatt M, Kress W, Williamson SJ, Nasko DJ, et al.
392 Ribonucleotide reductases reveal novel viral diversity and predict biological and
393 ecological features of unknown marine viruses. *Proc. Natl. Acad. Sci. National Academy*

- 394 of Sciences; 2014;111:15786–91.
- 395 29. McAllister SM, Davis RE, McBeth JM, Tebo BM, Emerson D, Moyer CL. Biodiversity
396 and emerging biogeography of the neutrophilic iron-oxidizing Zetaproteobacteria. *Appl.*
397 *Environ. Microbiol.* American Society for Microbiology (ASM); 2011;77:5445–57.
- 398 30. Biodocker. <http://biodocker.org/>. Accessed 2016 Jul 21.
- 399 31. Merkel D. Docker: lightweight Linux containers for consistent development and
400 deployment. *Linux J.* Belltown Media; 2014;2014:2.
- 401 32. GNU Operating System. <http://www.gnu.org/licenses/>. Accessed 2016 Jul 21.
- 402

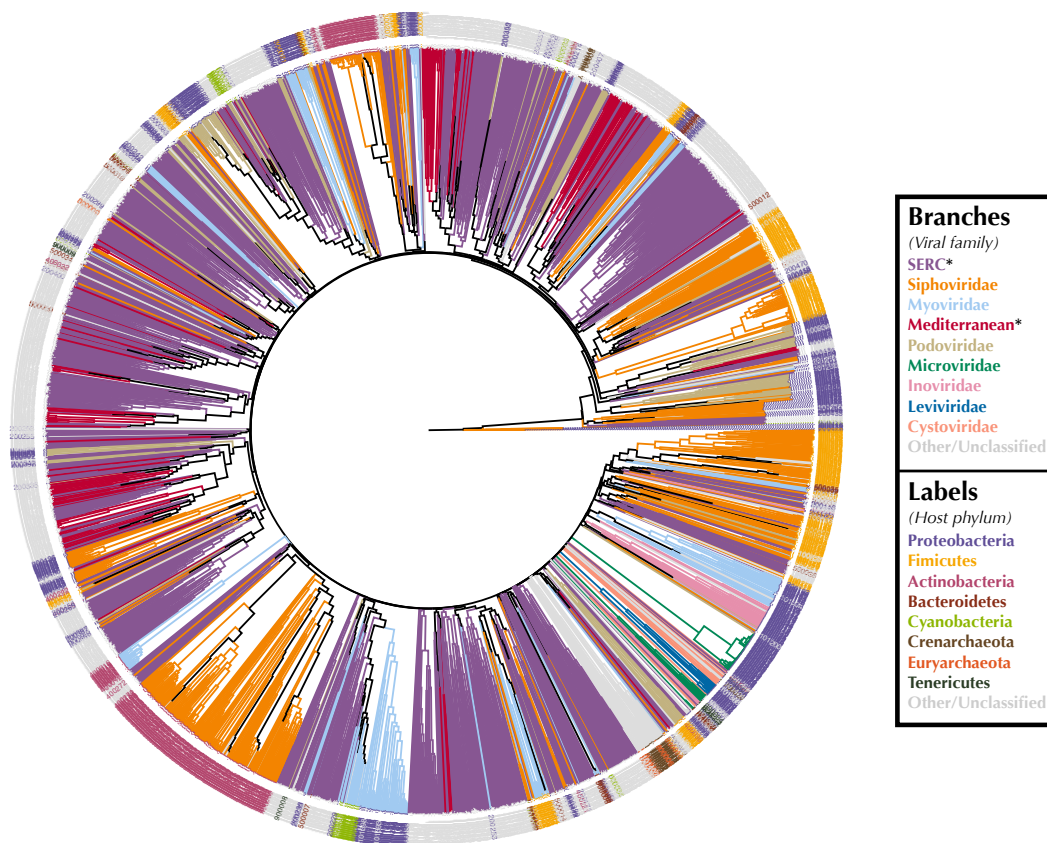


403

404 **Fig. 1: Precedence of Iroki's customization pipeline**

405 Flowchart illustrating the precedence of steps when performing multiple customizations
406 with Iroki. Input/output files are purple, command line options are in blue, and processes
407 are orange. The choice of exact or regular expression matching guides each subsequent
408 step of the process. Iroki gives higher precedence to processes towards the bottom of the
409 diagram. For example, given that a user selects the options for coloring both labels and
410 branches, and provides both a biom file and color map with the color map specifying
411 colors for the labels only, then the branches will be colored according to the color
412 gradient inferred from the biom file, whereas the labels will be colored according to the
413 rules specified in the color map.

414



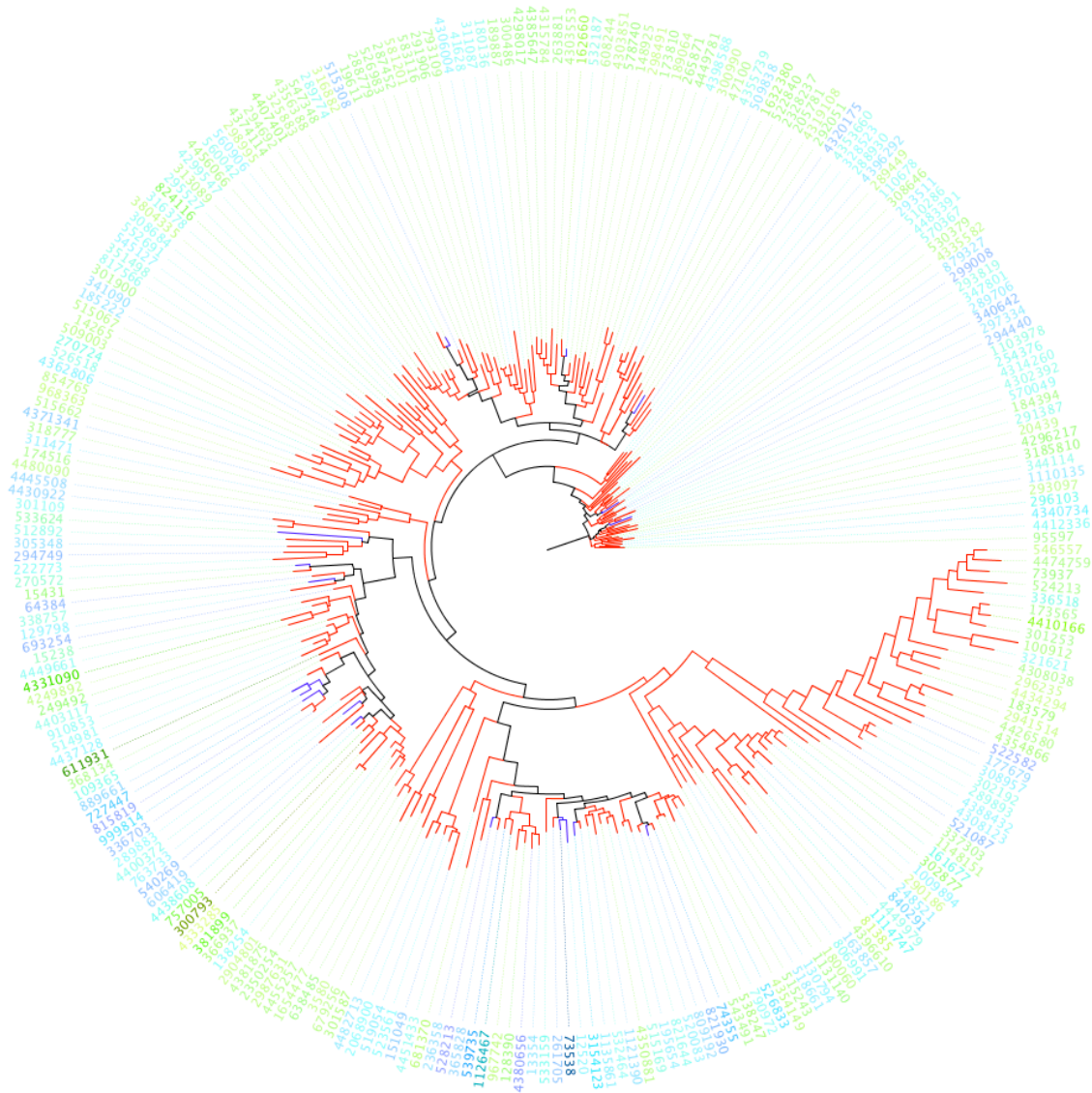
415

0.9

416 **Fig. 2: Comparing phage and their host phyla**

417 All phage genomes from Phage SEED with assembled virome contigs from the Chesapeake
418 Bay and Mediterranean Sea. Iroki highlights phylogenetic trends after coloring branches
419 according to viral family or sampling location in the case of virome contigs (marked with
420 an asterisk in the legend), and coloring node labels according to host phylum of the
421 phage.

422



423

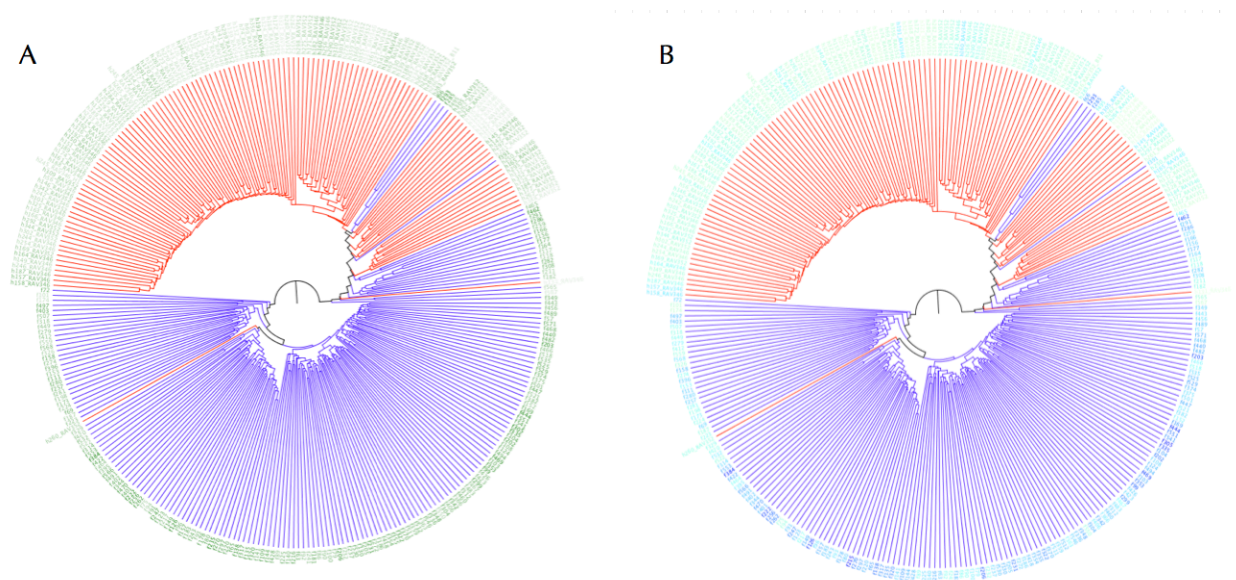
0.1

424 **Fig. 3: Changes in OTU abundance in two sample groups**

425 Approximate-maximum likelihood tree of OTUs that showed significant differences in
426 relative abundance between STEC positive and STEC negative cattle hide samples.

427 Branches show significance based on coloring by the p-value of a Mann-Whitney U test
428 examining changes in abundance between samples positive for STEC ($p < 0.05$ – red) and
429 samples negative for STEC, ($p \geq 0.05$ – blue). Label color on a blue-green gradient

430 highlights OTU occurrence based on the abundance ratio between STEC positive samples
431 (blue) and STEC negative samples (green). For example, labels that are darker green had a
432 higher abundance in STEC negative samples. Node luminosity represents overall
433 abundance with lighter nodes being less abundant than darker nodes.



434

435 **Fig. 4: Comparing cattle fecal and hide samples and the abundance of Ruminococcaceae**

436 Phylogeny based on UPGMA tree of pairwise unweighted UniFrac distance between 356

437 bacterial community profiles based on SSU rRNA amplicon sequences from cattle hide

438 and feces. Branches are colored by feces (blue) and hide (red). Rapid testing of the

439 hypothesis that the abundance of one of the most abundant families, Ruminococcaceae,

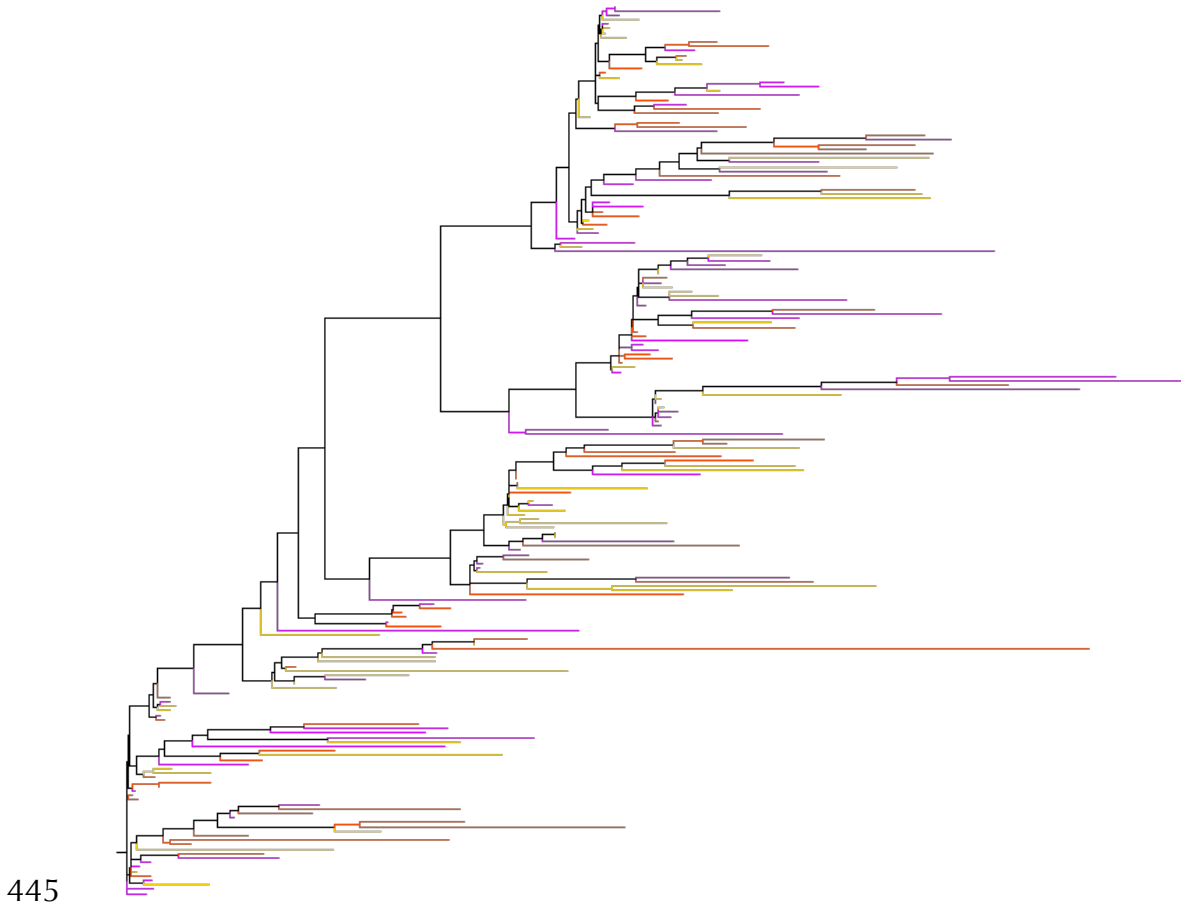
440 and sample origin are correlated is enabled through node label coloring by (A) a green

441 single-color gradient (color saturation increases with increasing abundance of

442 Ruminococcaceae OTUs) and (B) a light green (low abundance of Ruminococcaceae

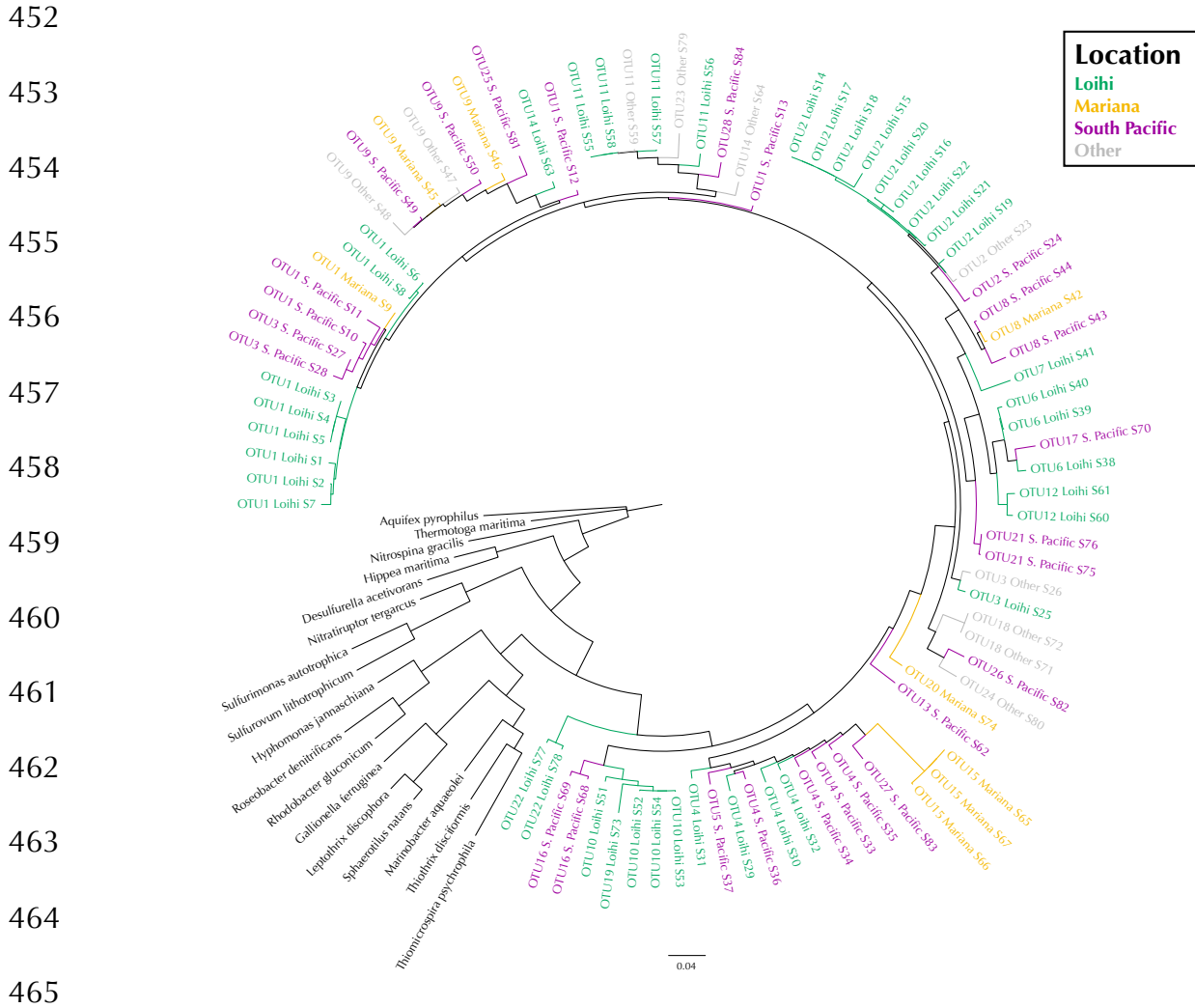
443 OTUs) to dark blue (high abundance of Ruminococcaceae OTUs) color gradient.

444



446 **Fig. 5: Temporal dynamics of virioplankton populations according to Cyano II RNR**
447 **amplicon phylogeny**

448 An approximately-maximum-likelihood phylogenetic tree of 200 randomly selected class
449 II Cyano RNR representative sequences from 98% percent clusters. Iroki was used to color
450 branches by time point: zero hours – yellow, six hours – orange, and twelve hours –
451 purple.



466 **Fig. 6: Zetaproteobacteria show biogeographic partitioning**

467 Phylogenetic tree showing placement of full-length Zetaproteobacteria SSU rRNA
468 sequences with outgroups. Iroki was used to color labels and branches by geographic
469 location of the sampling site (Loihi – green, Mariana – gold, South Pacific – purple, and
470 Other – gray), as well as to rename the nodes with OTU and sampling site metadata.

471